

ANNOUNCEMENTS —

1. HW6 due next Monday, 5pm / late due, next Wed, 5pm
 2. HW3 grades are (re)posted
-

Random thought: When do you think "group" exams would succeed?

- (a) if everyone in the group was an expert in ONE topic
- (b) more of a "jack-of-all-trades" situation?

(Mso, kehlamin for slope day?)

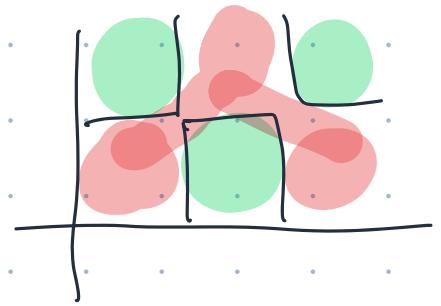
Recap: Decision trees

Goal: Have a non-linear classifier, akin to region-splitting

$J(R)$ measures the impurity in a region,

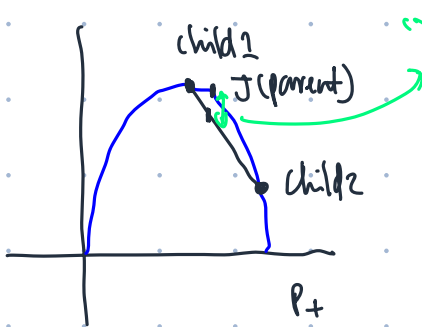
our goal:

$$\max J(\text{parent}) - \text{w. avg } J(\text{children})$$



yes/no questions

Entropy metric - measures impurity



parent entropy always above
w. avg of child entropies

⇒ splitting by entropy always results in
lower entropy.

$$J(R) = \sum_{y \in \{+1, -1\}} -p_y \log p_y$$

Alternative to not computing log, Gini impurity:

$$J_{\text{gini}}(R) = \sum_{y \in \{+1, -1\}} p_y (1 - p_y)$$

is latitude > 30°
/ \
 |

OR SPECIFIC INSTANCES

Today: ① variants of decision trees

② problems with decision trees

③ (As usual,) ask if we can do better?

Categorical attributes

If we split on netid, we'd be asking 399 questions. Data:

In the worst case, I'd be asking 2^{400} questions!



splitting on "netid" results in pure nodes
⇒ max decrease in "J"

Classify CS7780/5780 students into grad/ugrad classes

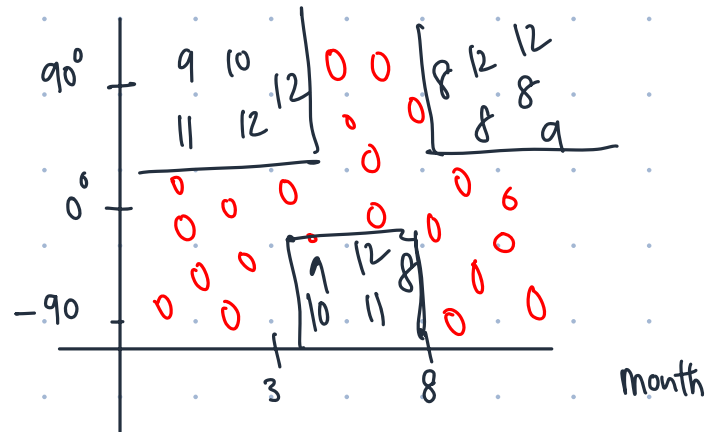
NetID
TF352
KS999
W692
LJ2
⋮

Takeaway: ⊕ Decision trees kind of give you nice way of handling categorical attributes

⊖ highly branching attributes are not favorable → computation
↳ overfitting.

Regression trees

Instead of predicting sled/not, we wish to predict snowfall in inches.



① Instead of a majority vote, we will use the avg in the region

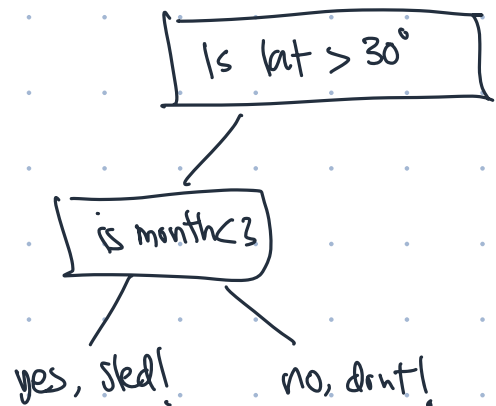
$$\hat{y} = \frac{1}{|R|} \sum_{j \in R} y_j^{(i)}$$

② "J" $\rightarrow J_{sav} = \frac{1}{|R|} \sum_{j \in R} (y_j^{(i)} - \hat{y})^2$

Why does CS3780/5780 like DT?

Rakia — They're non-linear!

Jay — Interpretability (for free!!)



Issues with Decision trees — Bias/Variance

Q. Do DT suffer from high bias? } Stopping criteria is that every leaf is a single, pure node.

⇒ I can ask questions (as many as I want) to get EVERY training pt. correctly classified!

⇒ **LOW BIAS!** ← good ;)

Q. Do DT suffer from high variance?

overfitting to the training data ⇒ **HIGH VARIANCE!**

BAD! ;(

⇒ Q. How do we fix this "high variance" problem?

(1) Set a min leaf size — don't split if $|R| < 20$

(2) set a max depth — if tree has > 3 levels, don't grow!

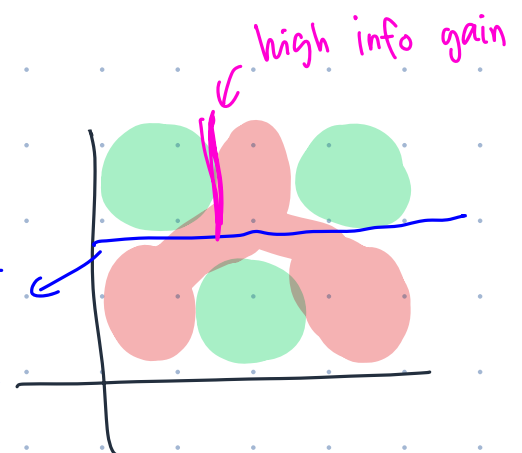
(3) threshold on max nodes — 20 nodes/less in the whole tree

(4) threshold on "info gain" we talked about

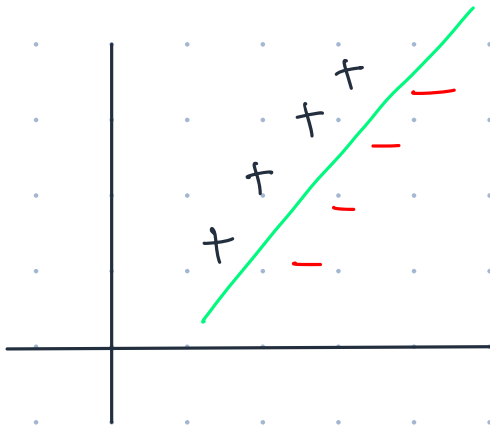
Don't

Jay: (5) grow out the whole tree, then remove nodes based on XYZ on validation set!

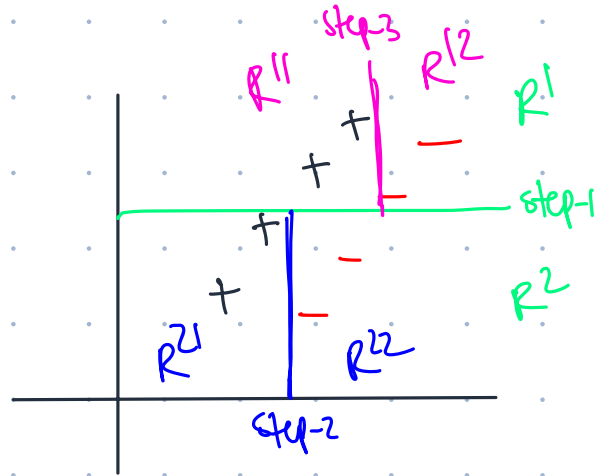
doesn't get me high info gain



"line drawing"



logistic regression



decision tree

Takeaway: Can't model additive structures

The good, the bad of DECISION TREES

In the wild, DT have poor predictive accuracy

Idea: what if we could train multiple models and use the "aggregate" prediction?

We have "m" total IID random variables, $x^{(j)}$ s,

we want to compute

$$\text{Var} \left[\frac{1}{m} \sum_{j=1}^m x^{(j)} \right] = \frac{\sigma^2}{m}, \text{ because } x^{(j)}\text{s are independent}$$

goal - variance reduction

idea - as $m \uparrow$ (# predictors) $\text{var} \downarrow (Y_m)$

} $x^{(j)}$ = error rate of my predictor $h^{(j)}$

Is the assumption of independence among predictors reasonable?

→ maybe not always

~~BREAK THE ASSUMPTION~~, $\text{Corr}(x^{(j)}, x^{(k)}) = \rho$

IID $x^{(j)}$ s (not necessarily ind),

$$\text{Var} \left[\frac{1}{m} \sum_{j=1}^m x^{(j)} \right] = \sigma^2 \rho + (1-\rho) \frac{\sigma^2}{m}$$

$\rho = 0$ when decorrelated

$\rho = 1$ we have σ^2

$\frac{\sigma^2}{m}$

ways to ensemble?

Ensemble "m" predictors &
hope var goes down!

1) Train different algorithms - SVM, DT, LR, NB - aggregate them!

Computationally expensive!

2) Go and collect different datasets - D^1, D^2, D^3 - train SVM on D^1, D^2, D^3
separately, aggregate

could be infeasible

similar to this, but doesn't require additional data collection! } "bagging"
(Eg. RF)

ensembling from bias-reduction
perspective = "boosting"
(eg. Adaboost, XGBoost)

BOOTSTRAPPING + AGGREGATION (AKA BAGGING)

Idea - If we had $D^1, D^2, D^3, \dots, D^m$, we could train a model on each of these separately, then aggregate!

$$D = \{(x^{(j)}, y^{(j)}) \mid 1 \leq j \leq n\}, \quad (x^{(j)}, y^{(j)}) \sim \mathcal{P} \leftarrow \text{true distribution}$$

1. Bootstrapping: We assume $D = \mathcal{P} \Rightarrow$ we can sample from \mathcal{D} , more imp, we can sample as many times as we want.

\Rightarrow Sample, WITH REPLACEMENT, 'n' samples, to get $z^{(1)}$
math - if we are assuming $D = \mathcal{P}$, sampling w/ repl makes sense for the assumpt to hold!

Repeat to get $z^{(2)}, z^{(3)}, \dots, z^{(m)}$

2. Aggregation: Given $z^{(1)}, \dots, z^{(m)}$, train some model on each, to get $h^{(1)}, \dots, h^{(m)}$

aggregate hypothesis \swarrow

$$h(x) = \frac{1}{m} \sum_{j=1}^m h^{(j)}(x)$$

Jay →
Niceties of BAGGING!

(1) Variance of average error of "m" correlated models,

$$\text{Var} \left[\frac{1}{m} \sum_{j=1}^m X^{(j)} \right] = \sigma^2 p + (1-p) \frac{\sigma^2}{m}$$

Thought - If we trained DT on $D, D, D \dots D$ m times
then "p" b/w $h^{(1)}, h^{(2)}, \dots, h^{(m)}$ would be high
corr ←
But, for bootstrap samples, p is lower!

If we sample a lot of $Z^{(j)}$ s, the $m \uparrow$, makes $(1-p) \frac{\sigma^2}{m} \downarrow$
↓
we are reducing the variance by training on a subset of the data,
bias \uparrow .

(2) Out-of-bag error.
In sampling $Z^{(j)}$, 33% of D doesn't get sampled!
free 'val' set

Almost

Random Forests

= DT + Bagging

DT + Bagging

= ALMOST RF

making DT ←
good choice to
bag with!

DT ← low bias, high variance

Bagging - increases bias,
decreases variance!

In DT construction, one of the attributes could be really important,
meaning all $h^{(1)}, \dots, h^{(m)}$ would likely split on
that feature (eg. HP Pokemon)



At the first step, all DT are correlated!

→ At each split, I am only going to use a random subset of
features



drives P , further down, since we're cutting
down on features,
we ↑ bias!

The good and the bad (of BAYESIAN)

Additive modeling is still issue!